# A NEW HOME FOR CALABASH

In summer 2020, Digital Scholarship Services at NYU Libraries was approached by NYU professor Jacqueline Bishop about finding a new home for **Calabash: A Journal of Caribbean Arts and Letters**. Multilingual and focused on centering unheard voices, Calabash was a pioneering journal showcasing poetry, literature, and visual arts from across the Caribbean. The journal, which Dr. Bishop edited from 2000–2008, had since ceased publishing, and the NYU server that had been hosting the site was due to be retired.

## NYU | LIBRARIES

### PROJECT TEAM

**Zach Coble** // Head, Digital Scholarship Services
**Jonathan Greenberg** // Digital Scholarly Publishing Specialist
**Marii Nyrop** // Senior Research Data Engineer
**Kate Pechekhonova** //Senior DevOps Engineer
**Alexandra Provo** // Research Curation Librarian
**Nick Wolf** // Research Data Management Librarian and Co-Head of Data Services
**Deb Verhoff** // Digital Collections Manager
**Vita Kurland** // NYU/LIU Dual Degree student
**Katherine Santana** // NYU/LIU Dual Degree student

### LINKS & CONTACT

**Alexandra Provo** | alexandra.provo@nyu.edu
**Calabash's new home** | http://hdl.handle.net/2451/62241
**Blog posts**
https://wp.nyu.edu/library-dlts/2022/09/14/a-new-home-for-calabash/
https://marii.info/notes/ghostscript-for-washed-out-pdfs

## 01 OPTIONS

Newly formed digital library governance groups were consulted and the following options proposed:

1. **Create a static site** (hosted by faculty member or by library)
2. **Extract PDFs and publish to the Faculty Digital Archive** (institutional repository built on DSpace)
3. **Build a new site** that would allow for better citation and discovery at the article level
4. **Journal is acquired by University Archives:** crawl using UA's instance of Archive-It, or treat as a born-digital archival collection.

## 02 DECISION

**Extract PDFs and deposit to institutional repository**

In contrast to digital editions with rolling releases, as a journal Calabash's content was released periodically: the website consisted of a landing page for each issue with links to PDFs. Although the creation of static sites is becoming a standard practice (see Endings Principles, 2023 and Ronallo, 2017), the focus for Dr. Bishop was on the articles, not necessarily the look and feel of the website. Another benefit of depositing in the institutional repository was the level of preservation commitment and improved citability via metadata, DOIs, and handle URLs.

## 03 INTERNET ARCHIVE

Although look and feel was not a priority, the presentation of a digital humanities project is part of its essence (Holmes and Takeda, 79). Since the site was already being crawled, we checked to see that the main pages were included in IA.

## 04 IMPLEMENTATION

### HTTrack WEBSITE COPIER

### Preparation

- Extracted PDFs from website
  - Copied off old server with httrack
  - OCR'd PDF documents and injected .txt back into PDFs
- Created a Metadata Application Profile (MAP) to configure DSpace collection

### Metadata acquisition and transformation

- Scraped website using Python
- Used OpenRefine and Google Sheets to combine spreadsheets, transform data, and derive filenames, volume and issue numbers
- Filled out DOI spreadsheets based on scraped metadata and ran a script to assign DOIs from our internal registry
  - Did manual quality check, comparing number of articles on website with number of rows in corresponding spreadsheet
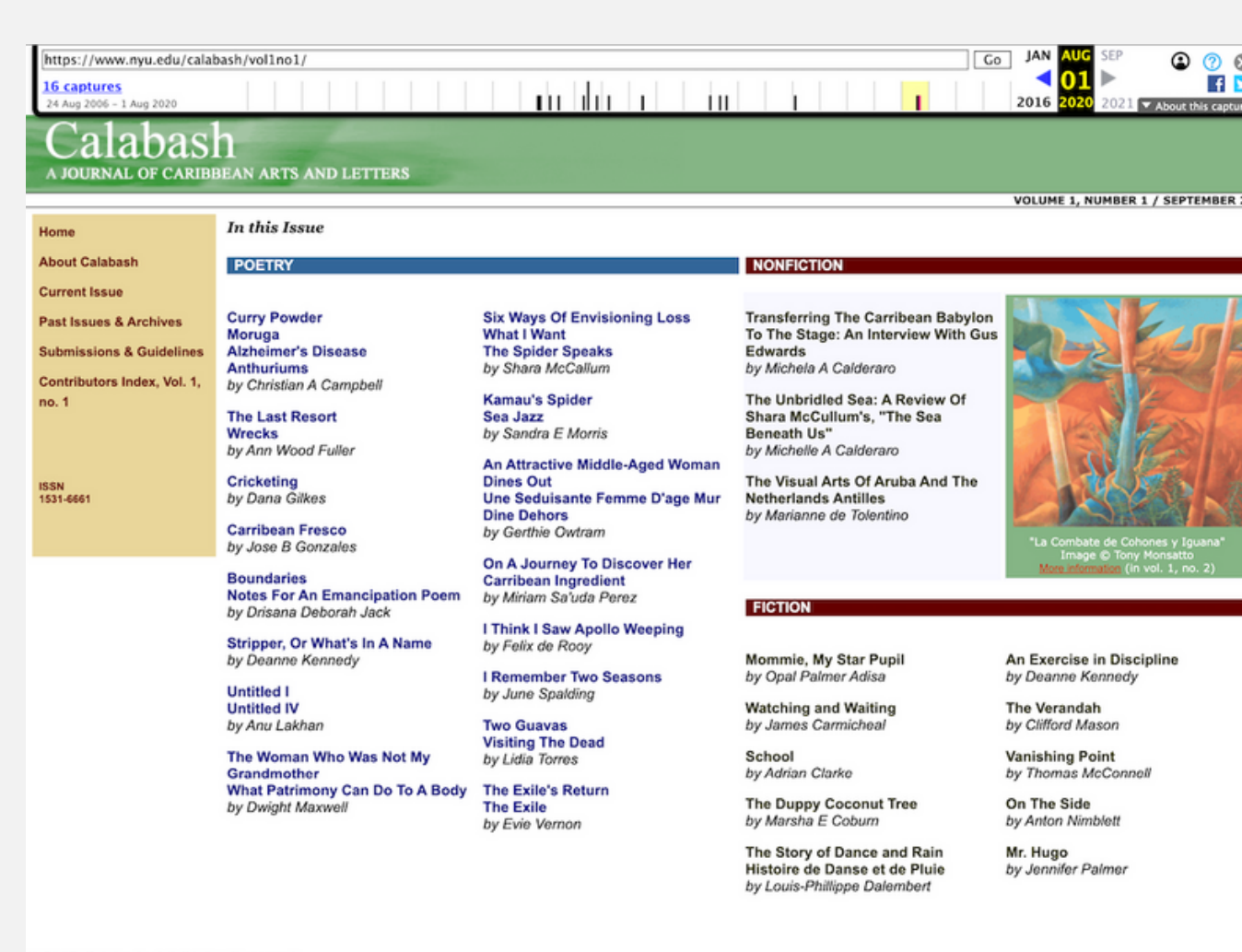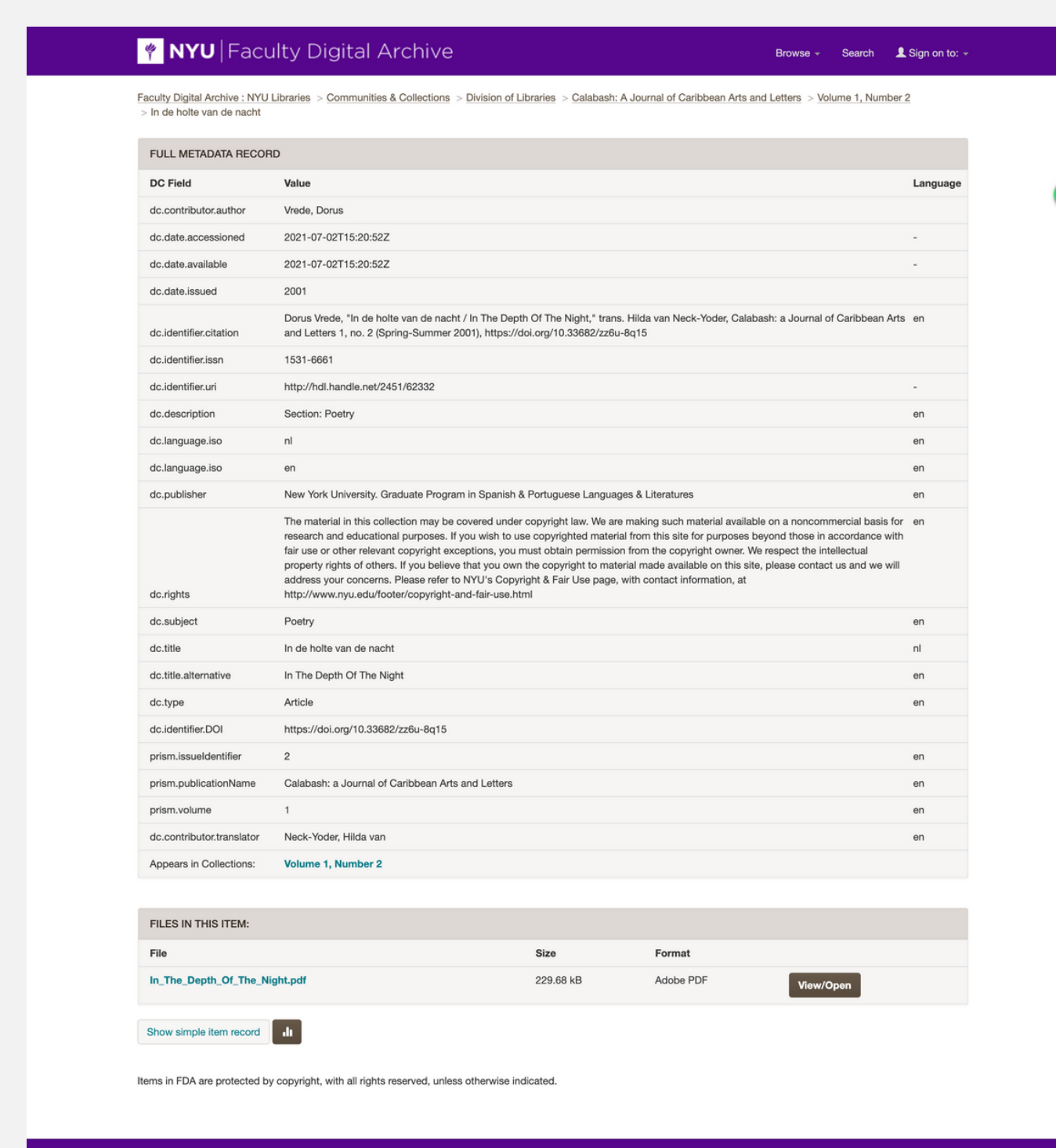
### Loading PDFs to FDA

- Added metadata spreadsheets to folders containing the OCR'd PDFs
- Used SAFBuilder to package XML metadata and bitstreams so we could upload both PDFs and metadata

### Metadata enrichment

- Two NYU/LIU Dual Degree students used the MAP to add information such as alternate titles, translators, subjects, and descriptions
- Focused on language, taking care to tag titles and alternate labels with their respective languages, and recording languages other than English in the language field

### Registering DOIs

- Generated Crossref XML files for DOI registration using Python scripts created in-house





## 05 REFLECTIONS

- **Webscraping**
  - Older, well-structured HTML website made for effective scraping
  - Typos were perpetuated, and not all metadata could be derived
- **Multilingual content**
  - DSpace uses ISO 639-1 (no language code for Papiamento)
  - Lack of language expertise on team
- **Temporally-extended and cross-departmental work is challenging**
  - Gap between journal's publication and preservation efforts
  - Departments with varied priorities means work happens in fits and starts

### Related Literature

"Endings Principles for Digital Longevity (Version 2.2.1)," March 3, 2023. https://endings.uvic.ca/principles/.
Holmes, Martin, and Joey Takeda. "From Tamagotchis to Pet Rocks: On Learning to Love Simplicity through the Endings Principles." Digital Humanities Quarterly 017, no. 1 (May 26, 2023).
Ronallo, Jason. "Sunsetting DH Projects." DLF Forum. Pittsburgh, PA, 2017. https://ronallo.com/presentations/sunsetting-dlf/slides-single-page.html.